# Citation Network Analysis

Thom Neale
Twitter: @twneale
Github: twneale

# Where to find materials

- The materials for this talk are all located at twneale.github.io/citation-network-analysis

./ **Citation Network Analysis**
Materials for my talks at PyData Boston 2013 and Law Via the Internet 2013

Download as .zip    Download as .tar.gz    View on GitHub

This page has links to my talk materials and other resources, contact info, and further reading.

My paper on Citation Analysis of Canadian Case Law.

An IPython Notebook with runnable network analysis code.

Reading

>>   Programming the Semantic Web
>>   Graph Databases

Contact

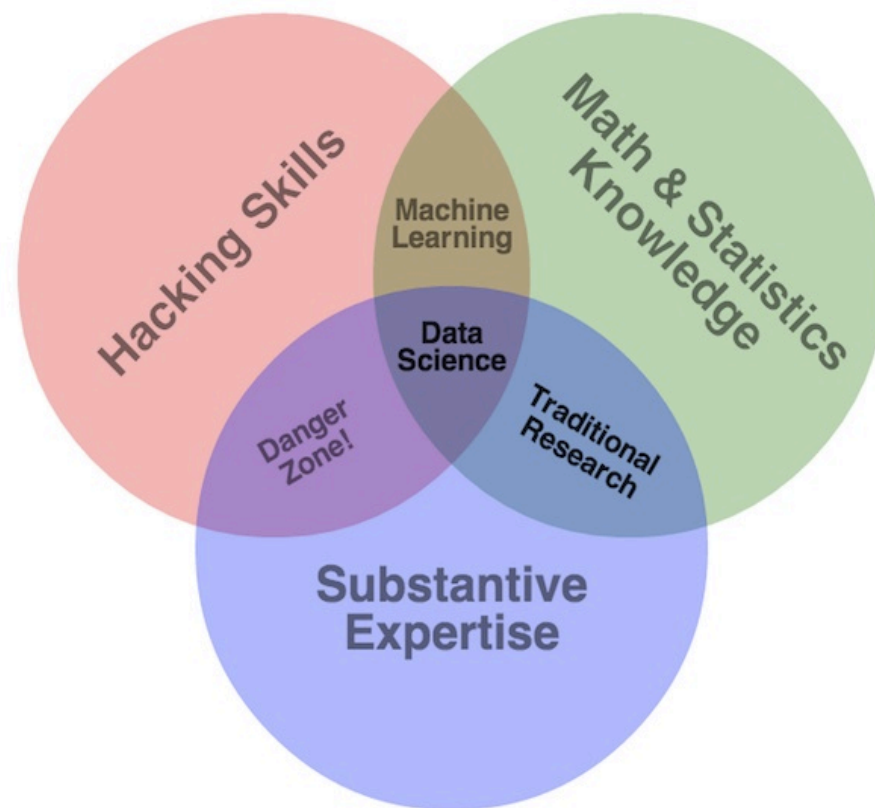>>   Twitter: @twneale
>>   Github: twneale

# Goodies

- Webpage where you can try out the code I used in this study: https://www.wakari.io/twneale

- App demonstrating (sort of) practical use of network analysis data: https://cite-fight.com

# Big Picture

- How to get from unstructured text to data that people use to create excellent tools

Hacking Skills

Math & Statistics Knowledge

Machine Learning

Data Science

Danger Zone!

Traditional Research

Substantive Expertise

drewconway.com

# This talk is about

- The journey from unstructured text to highly structured data

- Fantasizing about the amazing things we can do once the journey is complete

Consider the following question:

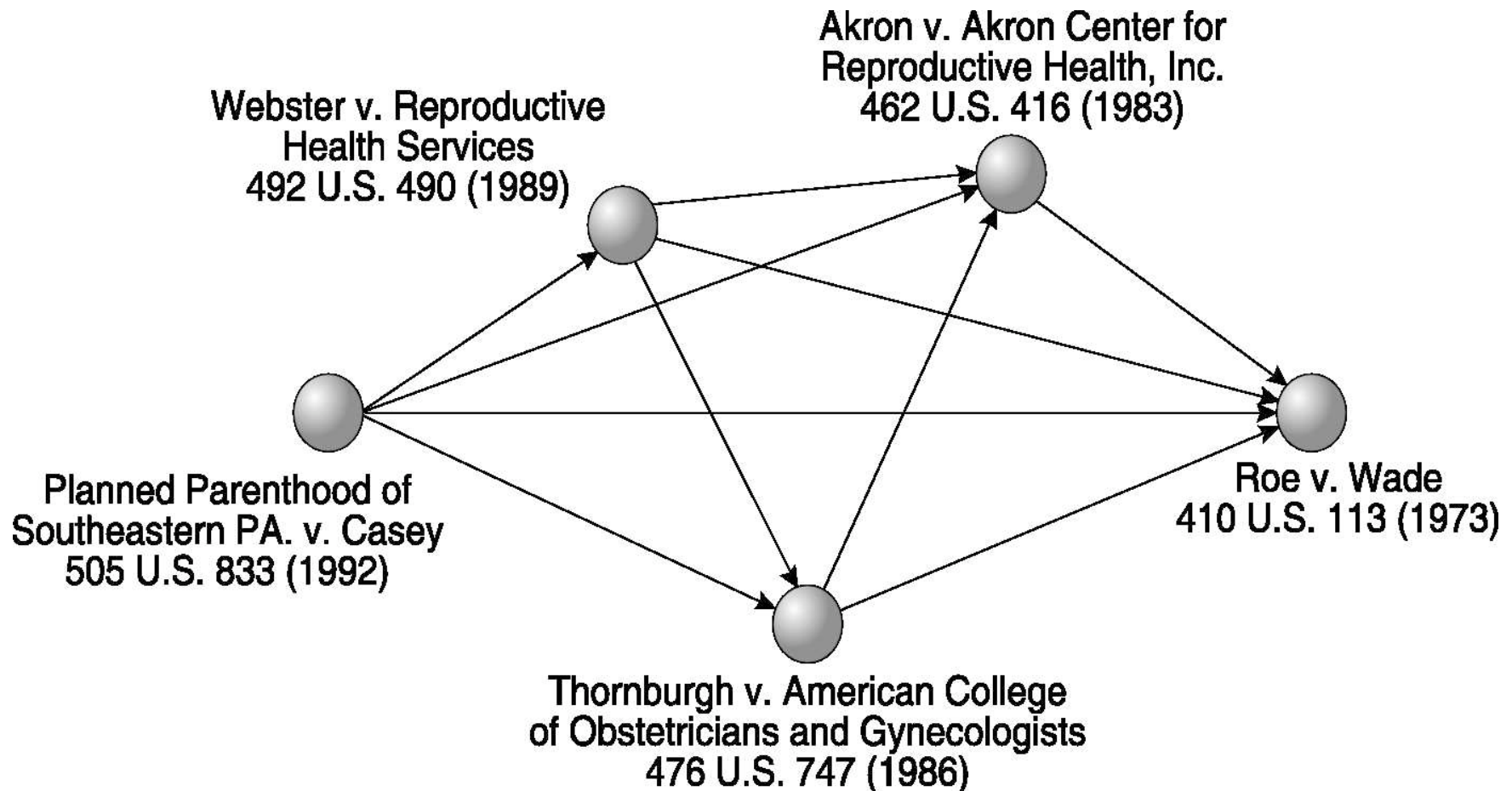# What determines the importance of a case?

# Maybe this:

Michael Gerhardzt observed that the extent and nature of a precedent's network of citations determine the strength of its constraining power on subsequent cases. He argued further that the authority of a precedent depends on the consistency and uniformity with which other authorities have cited it.

Michael J. Gerhardt, The Irrepressibility of Precedent, 86 N.C.L.REV.1279, 1291 (2008)
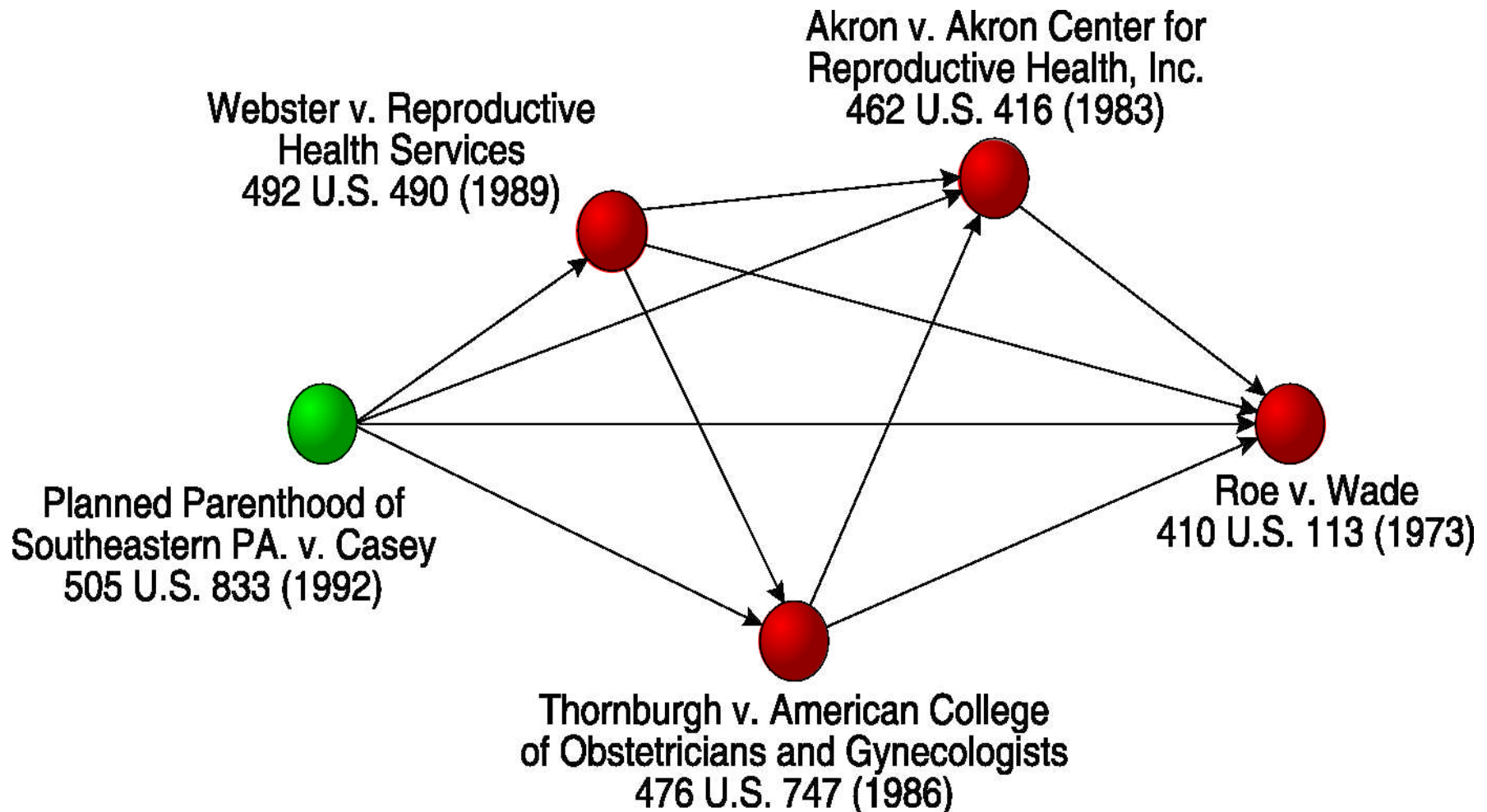
# Network of citations?



Akron v. Akron Center for
Reproductive Health, Inc.
462 U.S. 416 (1983)

Webster v. Reproductive
Health Services
492 U.S. 490 (1989)

Planned Parenthood of
Southeastern PA. v. Casey
505 U.S. 833 (1992)

Roe v. Wade
410 U.S. 113 (1973)

Thornburgh v. American College
of Obstetricians and Gynecologists
476 U.S. 747 (1986)

Fowler et at, *The Authority of Supreme Court Precedent* (2008)

# Translation:

- Constraining power of a case on subsequent cases (i.e., importance) depends on two things:

1) Nature and extent of case's network of citations

2) Consistency and uniformity with which other authorities have cited it

# 1) The extent and nature of a precedent's network of citations



Akron v. Akron Center for
Reproductive Health, Inc.
462 U.S. 416 (1983)

Webster v. Reproductive
Health Services
492 U.S. 490 (1989)

Planned Parenthood of
Southeastern PA. v. Casey
505 U.S. 833 (1992)

Roe v. Wade
410 U.S. 113 (1973)

Thornburgh v. American College
of Obstetricians and Gynecologists
476 U.S. 747 (1986)

# Conclusion

*Planned Parenthood v Casey* is
"well founded in law"

# 2) Consistency and uniformity with which other authorities have cited it.



Akron v. Akron Center for
Reproductive Health, Inc.
462 U.S. 416 (1983)

Webster v. Reproductive
Health Services
492 U.S. 490 (1989)

Planned Parenthood of
Southeastern PA. v. Casey
505 U.S. 833 (1992)

Roe v. Wade
410 U.S. 113 (1973)

Thornburgh v. American College
of Obstetricians and Gynecologists
476 U.S. 747 (1986)

# Conclusion:

*Roe v Wade* is "influential"

# Turn that into an algorithm

# Make it recursive...

# Jon M. Kleinberg

Authoritative Sources in a Hyperlinked
  Environment (1998)

http://www.cs.cornell.edu/home/kleinber/auth.pdf

"Hyperlink-induced Topic Search (HITS)

# PageRank is another

# Degree Centrality is Another

# Recap

# Case Citations form a Network



Akron v. Akron Center for
Reproductive Health, Inc.
462 U.S. 416 (1983)

Webster v. Reproductive
Health Services
492 U.S. 490 (1989)

Planned Parenthood of
Southeastern PA. v. Casey
505 U.S. 833 (1992)

Roe v. Wade
410 U.S. 113 (1973)

Thornburgh v. American College
of Obstetricians and Gynecologists
476 U.S. 747 (1986)

# We can use network analysis algorithms to rank the nodes in the network



Ravellaw.com hotness

# Numerous Algorithms Exist

- Indegree Centrality

- PageRank

- HITS (Hyperlink-Induced Topic Search)

- Eigenvector Centrality

- And many, many, other stupefying algorithms

# The network of case citations is "scale-free"



Full Caselaw Network Degree Distribution

# Scale-free network



Airline*
— Emirates
— Lufthansa
— Delta Air Lines
— Air France
— British Airways
— Cathay Pacific
— Singapore Airlines

*Ordered by scheduled international passenger kilometres
flown in 2010 (source Wikipedia).
Only routes in the OpenFlights database are plotted.

Map: James Cheshire, spatialanalysis.co.uk
Flights Data: openflights.org
Basemap Data: naturalearthdata.com

# Random network

# Other important scale-free networks?

# THE INTERNET

# Network rankings change over time

# Time-series data creates interesting opportunities

# You can fit curves to the data



Citation: 1978canlii368 [1978]

Degree 1, r2=0.21619914545792074
Degree 2, r2=0.40066314267010467
Degree 3, r2=0.51759400395896082
Degree 4, r2=0.57952987287326108
Degree 5, r2=0.59482315868527968

# The curves enable you to estimate things



Citation: 1986canlii29 [1986]

Degree 1, r2=0.50375268306034449
Degree 2, r2=0.51875518347489746
Degree 3, r2=0.53242250587852358
Degree 4, r2=0.53840667786068763
Degree 5, r2=0.5391426445588744

Citation: 1975canlii146 [1975]

Degree 1, r2=0.067829170976794334
Degree 2, r2=0.51416432722385141
Degree 3, r2=0.6298901976220499
Degree 4, r2=0.65174652438747549
Degree 5, r2=0.69546728528452584

Citation: 1978canlii11 [1978]

Degree 1, r2=0.74647907454851548
Degree 2, r2=0.93817983489830148
Degree 3, r2=0.95003710165115562
Degree 4, r2=0.95256907875935137
Degree 5, r2=0.95728698713667881

Citation: 1987canlii84 [1987]

Degree 1, r2=0.71160214839066671
Degree 2, r2=0.71660525835354905
Degree 3, r2=0.80963009745276593
Degree 4, r2=0.93566052987475834
Degree 5, r2=0.93781899255667744

# What things?

- Is this case's importance increasing or decreasing? (slope, derivative)

- Which of these cases has had a greater cumulative influence over time? (area)

- Does anyone still use this case? (x intercept)

# Is this case's importance increasing or decreasing?

# Which of these cases has had a greater cumulative influence over time?

# Does anyone still use this case?



Citation: 1978canlii368 [1978]

Degree 1, r2=0.21619914545792074
Degree 2, r2=0.40066314267010467
Degree 3, r2=0.51759400395896082
Degree 4, r2=0.57952987287326108
Degree 5, r2=0.59482315868527968

# Findings

- On average, Canada Supreme Court (CSC) cases "fail" after 50 years.

- About 18% of CSC cases have survived longer than 15 years (and is still positive?).

- In all other courts, the average time to failure ranges from 3 to 15 years.

- In all other courts, less than 3% of cases survive longer than 15 years.

# Challenges

- Citation Extraction is hard
- Resolving citations to sources is harder

# Citation Extraction is hard

- Regular expressions are quick and easy
- But they don't scale
- Regexes alone aren't good at processing highly variable patterns
- "citation parsing" is a special case of entire document parsing
- Nested data structures (like citations) require stateful, recursive code

# Regexes are neither stateful nor recursive

# This

"There is some case law suggesting (without much discussion) that a purchaser cannot maintain a caveat unless it can be shown that specific performance is available. Where there is no binding contract, such that the purchaser is unable to get any remedy, clearly a caveat cannot be maintained: Oxford Development Group Inc. v. Midland Development Ltd., [1993] A.J. No. 47 (C.A.)."

# This

"There is some case law suggesting (without much discussion) that a purchaser cannot maintain a caveat unless it can be shown that specific performance is available. Where there is no binding contract, such that the purchaser is unable to get any remedy, clearly a caveat cannot be maintained: Oxford Development Group Inc. v. Midland Development Ltd., [1993] A.J. No. 47 (C.A.)."

# This

"There is some case law suggesting (without much discussion) that a purchaser cannot maintain a caveat unless it can be shown that specific performance is available. Where there is no binding contract, such that the purchaser is unable to get any remedy, clearly a caveat cannot be maintained: Oxford Development Group Inc. v. Midland Development Ltd., [1993] A.J. No. 47 (C.A.)."

# Becomes this

```
(0, Token.Content, u'There is . . .cannot be maintained: '),
(297, Token.Title, u'Oxford Development Group Inc. v. Midland Development Ltd.'),
(356, Token.SlipYear, u'[1993]'),
(363, Token.Reporter, u'A.J. No.'),
(372, Token.SlipNumber, u'47'),
(375, Token.ParenAbbrev, u'(C.A.)'),
(381, Token.Content, u'; '),
```

# Then This

```
-Node([])
  -Content([[(0, Token.Content, u'There is
some case law suggesting (without much discussion) that a purchaser
cannot maintain a caveat unless it can be shown that specific
performance is available. Where there is no binding contract, such
that the purchaser is unable to get any remedy, clearly a caveat
cannot be maintained: ')])
  -Source([])
    -Title([[(297, Token.Title, u'Oxford Development Group Inc. v.
Midland Development Ltd.')])
      -Citations([])
        -Citation([])
          -SlipYear([[(356, Token.SlipYear, u'[1993]')])
          -Reporter([[(363, Token.Reporter, u'A.J. No.')])
          -SlipNumber([[(372, Token.SlipNumber, u'47')])
          -Jurisdiction([[(375, Token.ParenAbbrev, u'(C.A.)')])
  -Content([[(381, Token.Content, u'; ')])
```

# Hey wait! That's a network!

# Citation Data Should Probably Be Stored as a Graph

- Complex queries will be cheap, rather than impossible

- More information gets stored, probably in less space

# Really Hard

- Resolving citations back to sources

- Why:

  - Volume pages are not unique identifiers

  - Titles aren't unique identifiers either

  - An alarming percentage of citations contain typos

# Volume Pages aren't unique IDs

# Comparing Titles

- Is unreliable
- Even hard for

# At Minimum

- To resolve book citations to sources, you need:

    - Detailed metadata about the book volumes

    - Need it PER VOLUME! Yes, it can change from volume to volume

    - Are cases tabular? Full text? Discretely paginated? Continuously paginated?
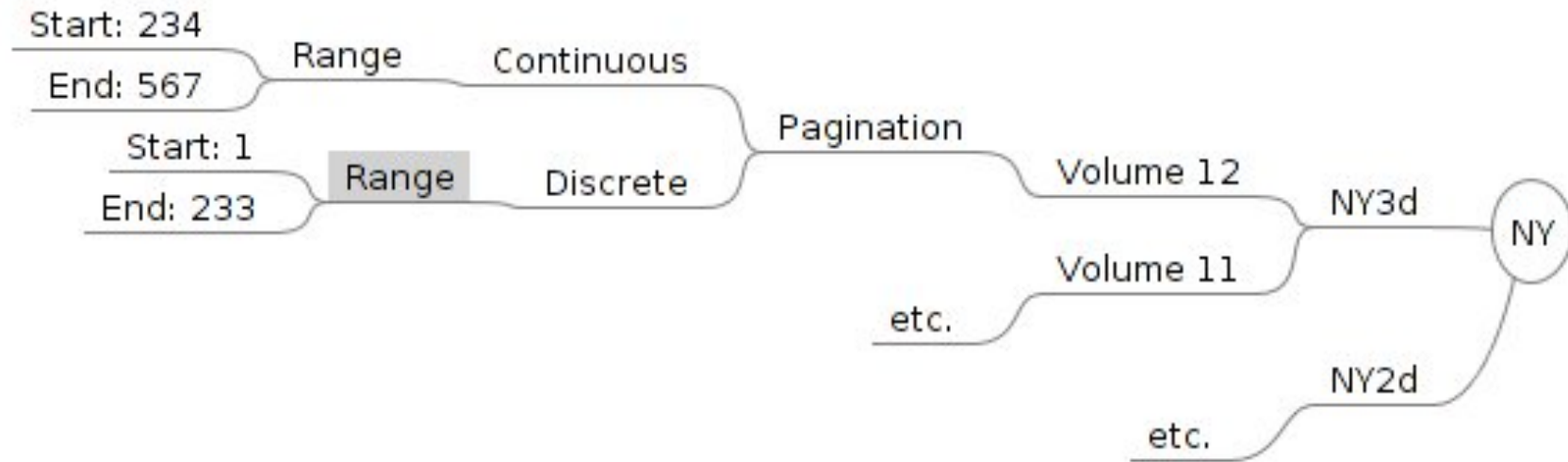
# Need Publication Metadata



```python
File  Edit  Selection  Find  View  Goto  Tools  Project  Preferences  Help

cosponsor_pagerank.py ×   effectiveness.py ×   demo.html ×   generate_json.py ●   events.py ×   bills.py — oh ×   __init__.py — oh ×   bills.py — nj ×   __init__.py — ca ×   bills.py — ma ×   votes.py ×        import datetime ●   __init__.py — al ×

34
35  REPORTERS = {'A.': [{'cite_type': 'state_regional',
36                       'editions': {'A.': (datetime.date(1885, 1, 1),
37                                           datetime.date(1938, 12, 31)),
38                                    'A.2d': (datetime.date(1938, 1, 1),
39                                             datetime.date(2010, 12, 31)),
40                                    'A.3d': (datetime.date(2010, 1, 1),
41                                             datetime.date.today())},
42                       'mlz_jurisdiction': 'us',
43                       'name': 'Atlantic Reporter',
44                       'variations': {'A. 2d': 'A.2d',
45                                      'A. 3d': 'A.3d',
46                                      'A.R.': 'A.',
47                                      'A.Rep.': 'A.',
48                                      'At.': 'A.',
49                                      'Atl.': 'A.',
50                                      'Atl.2d': 'A.2d',
51                                      'Atl.R.': 'A.'}}],
52               'A.D.': [{'cite_type': 'state',
53                         'editions': {'A.D.': (datetime.date(1896, 1, 1),
54                                               datetime.date(1955, 12, 31)),
55                                      'A.D.2d': (datetime.date(1955, 1, 1),
56                                                 datetime.date(2004, 12, 31)),
57                                      'A.D.3d': (datetime.date(2003, 1, 1),
58                                                 datetime.date.today())},
59                         'mlz_jurisdiction': 'us;ny',
60                         'name': 'New York Supreme Court Appellate Division Reports',
61                         'variations': {'A.D. 2d': 'A.D.2d',
62                                        'A.D. 3d': 'A.D.3d',
63                                        'AD 2d': 'A.D.2d',
64                                        'AD 3d': 'A.D.3d',
65                                        'Ap.': 'A.D.',
66                                        'Ap.2d.': 'A.D.',
67                                        'App.Div.': 'A.D.',

Line 34, Column 1                                                          Spaces: 2        Python
```
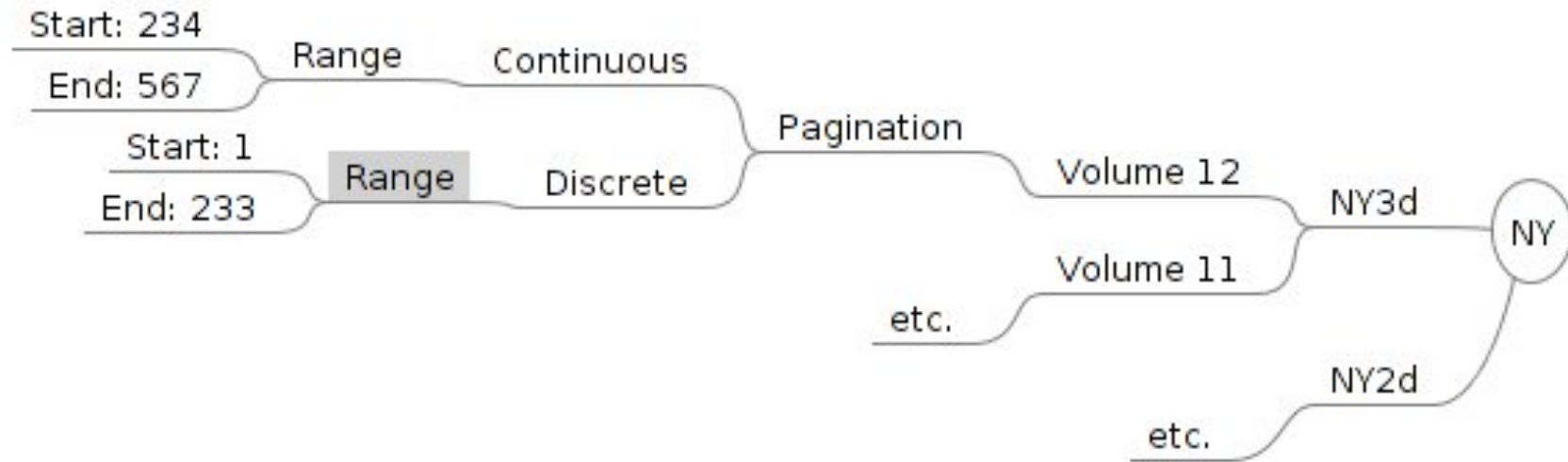
# It might look like this
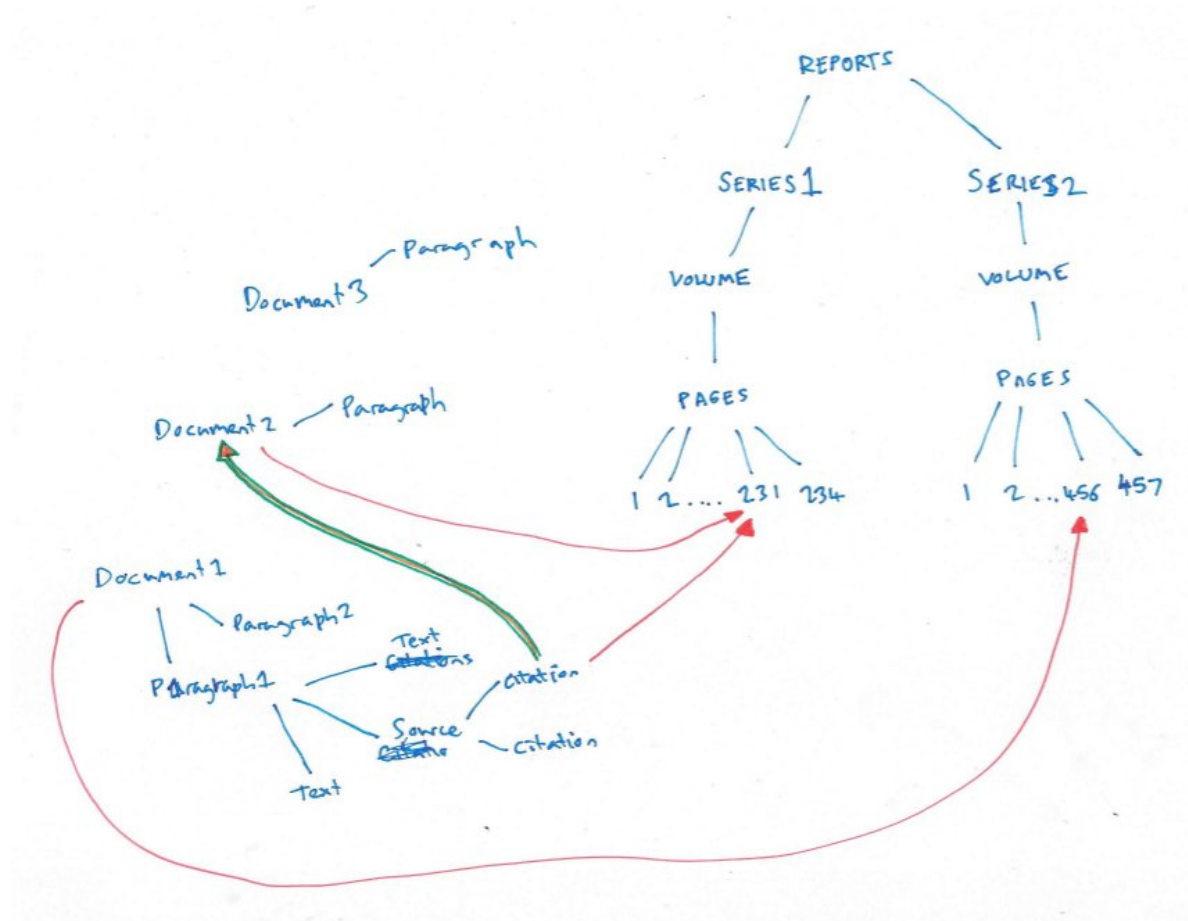
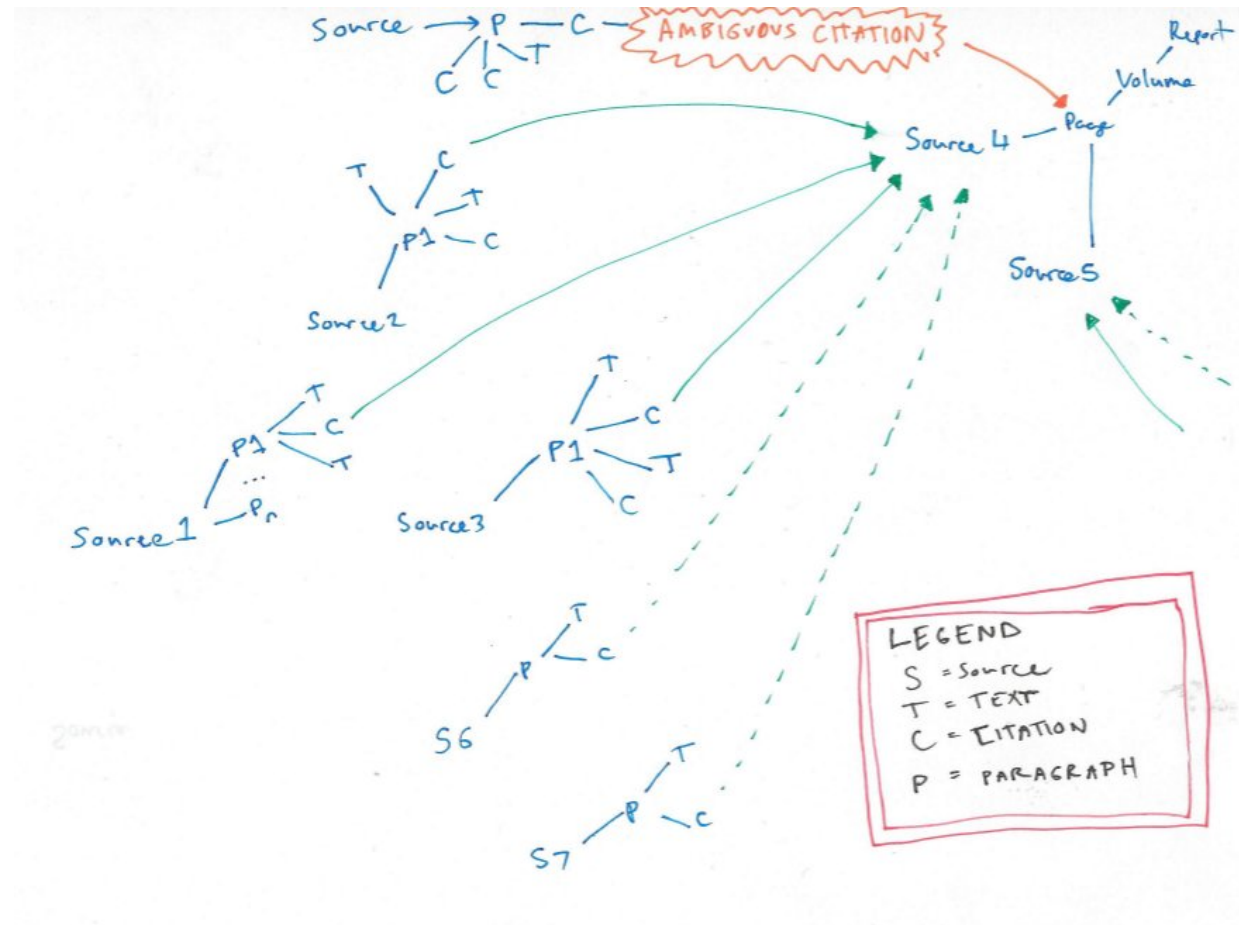# Wait a minute...

# It might look like this

# No Thom

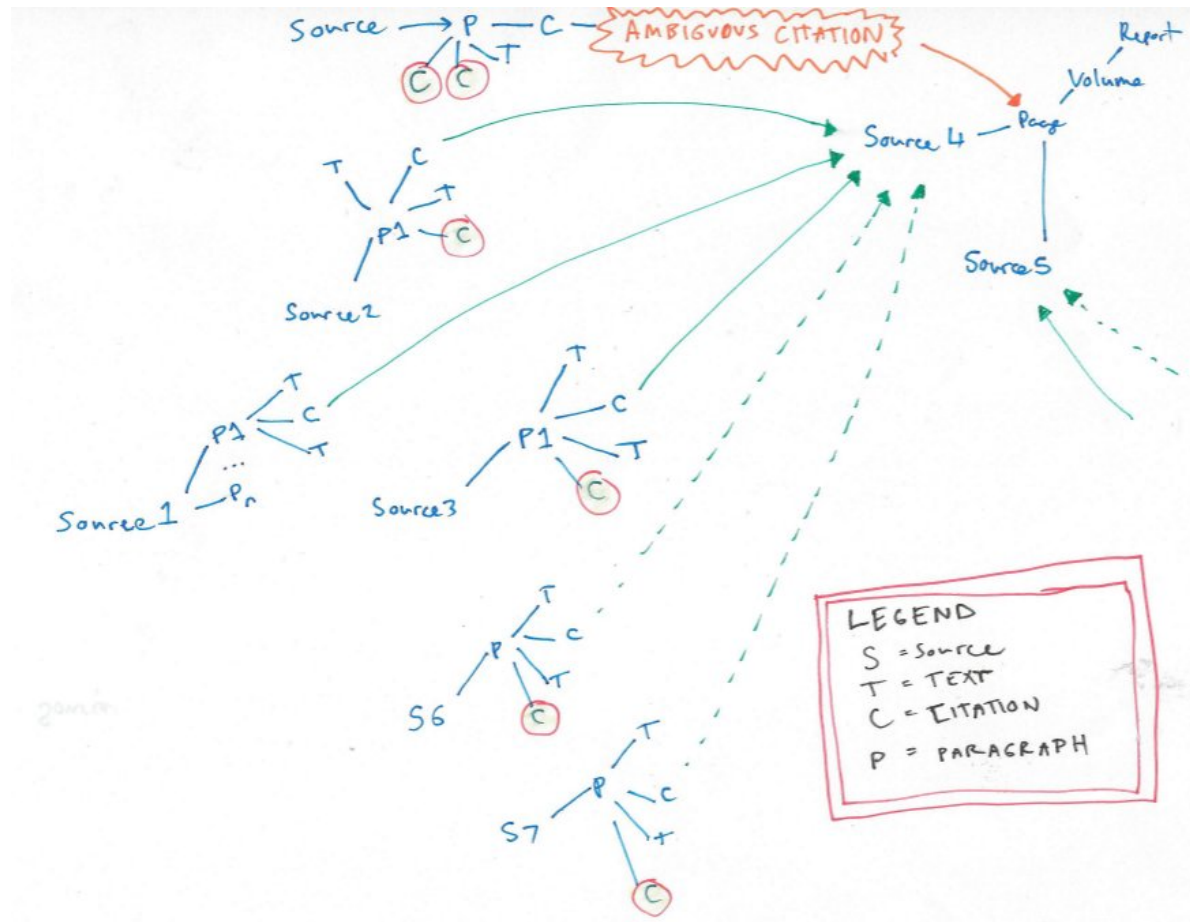None of this is useful. Take a yellow pad back to your office and create some value for once.

# An Overcomplicated Hypo

# Instead of comparing titles

# Ask the neighbors

# Recap

- We discussed a strategy for resolving ambiguous book citations that didn't require title comparison

- It was only possible because our graph database contains 1) publication metadata (report, series, volume, pagination), 2) cases, structured as subgraphs

# Not So Fast

- A horrifying percentage of citations contain typos, and resolving them to a source because even more ridiculous

- We have to try to reverse engineer the typo, then repeat the process for each candidate

- Title comparison would probably be helpful here too

# The Technology

- Python
- Networkx
- EC2
- Celery
- Numpy, SciPy
- neo4j

# Further Reading

- Programming the Semantic Web (Toby Segaran)

- Collective Intelligence (Toby Segaran)

- Graph Databases (Ian Robinson)

- Machine Learning for Hackers (Drew Conway)

# The End

# Citation Network Analysis

Thom Neale
Twitter: @twneale
Github: twneale

twneale.github.io/citation-network-analysis